

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-85710

(43) 公開日 平成11年(1999) 3月30日

(51) Int.Cl.⁶

G 0 6 F 15/16

識別記号

3 7 0

F I

G 0 6 F 15/16

3 7 0 M

審査請求 未請求 請求項の数 6 O L (全 12 頁)

(21) 出願番号 特願平9-250249

(22) 出願日 平成9年(1997) 9月16日

(71) 出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(72) 発明者 浅野 滋博

神奈川県川崎市幸区小向東芝町1番地 株

式会社東芝研究開発センター内

(72) 発明者 金井 達徳

神奈川県川崎市幸区小向東芝町1番地 株

式会社東芝研究開発センター内

(72) 発明者 菅野 伸一

神奈川県川崎市幸区小向東芝町1番地 株

式会社東芝研究開発センター内

(74) 代理人 弁理士 外川 英明

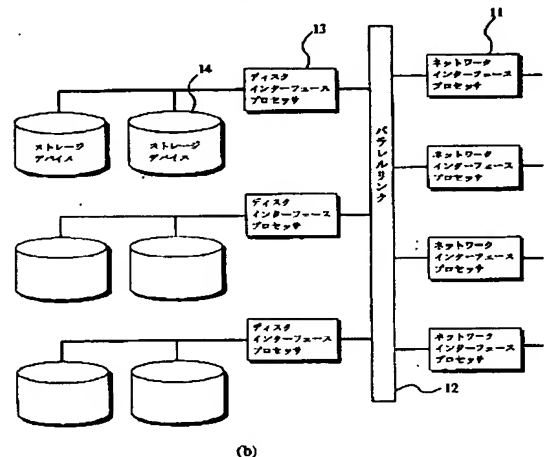
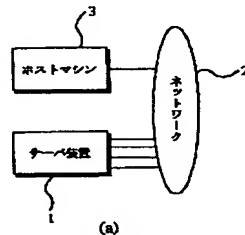
最終頁に続く

(54) 【発明の名称】 サーバ装置およびファイル管理方法

(57) 【要約】

【解決課題】 本発明は、管理が簡単でかつコストパフォーマンスに優れたサーバを提供することを目的とするとともに、階層キャッシュシステムにおける効率のよいファイル管理方法をも提供することを目的とする。

【解決手段】 本発明は、少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置(11)と、前記複数の第1の装置と転送手段(12)を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置(13)と、前記複数の第2の装置の各々に接続された第3の記憶手段(14)とを具備し、前記第1の装置に対して与えられるネットワークからの要求に基づく処理を、所定の条件に応じて、前記第1の記憶手段、前記第2の記憶手段または前記第3の記憶手段のうちのいずれかの記憶手段に対して行うことを特徴とするサーバ装置である。



【特許請求の範囲】

【請求項1】 少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備し、前記第1の装置に対して与えられるネットワークからの要求に基づく処理を、所定の条件に応じて、前記第1の記憶手段、前記第2の記憶手段または前記第3の記憶手段のうちのいずれかの記憶手段に対して行うことを特徴とするサーバ装置。

【請求項2】 ネットワークから要求を与えられる前記第1の装置は、該要求に関連する所定の情報に基づいて、前記第2の装置を選択することを特徴とする請求項1記載のサーバ装置。

【請求項3】 前記第3の記憶手段から前記第1の装置に所定のデータが送出される際に、所定の計算を行い、該得られた計算結果を該所定のデータとともに、前記第1の記憶手段に記憶することを特徴とする請求項1記載のサーバ装置。

【請求項4】 少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備し、前記第1の記憶手段に第1のデータが記憶され、前記第2の記憶手段に該第1のデータに対応する第2のデータが記憶されている場合に、前記第1の記憶手段から該第1のデータが追い出されるときは、前記第2の記憶手段に記憶された第2のデータの追い出し順位を他のデータの追い出し順位よりも低く設定するための処理を行うことを特徴とするサーバ装置。

【請求項5】 少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備するサーバ装置におけるファイル管理方法であって、前記第1の装置に対して与えられるネットワークからの要求に基づく処理を、所定の条件に応じて、前記第1の記憶手段、前記第2の記憶手段または前記第3の記憶手段のうちのいずれかの記憶手段に対して行うことを特徴とするファイル管理方法。

【請求項6】 少なくとも第1のプロセッサ及び第1の

記憶手段を有する複数の第1の装置と、

前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、

前記複数の第2の装置の各々に接続された第3の記憶手段とを具備するサーバ装置におけるファイル管理方法であって、

前記第1の記憶手段に第1のデータが記憶され、前記第2の記憶手段に該第1のデータに対応する第2のデータが記憶されている場合に、前記第1の記憶手段から該第1のデータが追い出されるときは、前記第2の記憶手段に記憶された第2のデータの追い出し順位を他のデータの追い出し順位よりも低く設定するための処理を行うことを特徴とするファイル管理方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、サーバ装置およびファイル管理方法に関するものである。

【0002】

【従来の技術】インターネットの発達により情報の流通の速度が飛躍的に高まり、ますます多くの人がインターネットにアクセスするようになってきている。情報を提供する手段としてはWorld Wide Web（以下「Web」という）によるものが一般的であり、プロトコルはHTTPおよびTCP/IPが使用される。情報を多くの人に提供するためには、ネットワークのバンド幅の拡大だけでなく、情報を蓄積したサーバがネットワークに送り出す能力の拡大も同時に必要である。

【0003】サーバに必要な機能は、主として、ストレージデバイスに蓄積された情報をネットワークに送り出す機能であるが、ストレージデバイスに蓄積された情報のフォーマットとネットワークに送り出す情報のフォーマットが異なるために、フォーマットの変換を行わなければならない。また、HTTPやTCP/IPで通信するためにプロトコルを実行する必要がある。

【0004】従って、これらの処理を1台のプロセッサで行っていたのではプロセッサにかかる負荷が増大し、性能向上のボトルネックとなる。また、前述したフォーマット変換のためにプロセッサがメモリ内で情報をコピーすることが頻繁に起こるので大きなメモリバンド幅が必要になる。

【0005】これらのボトルネックを緩和するためにプロセッサおよびメモリを複数結合するマルチプロセッサ方式のサーバが知られている。例えば、分散PC方式は、パフォーマンスを向上させるために、ネットワーク上に複数のWebサーバを設置したものである。しかしながら、この分散PC方式では、複数のWebサーバのそれぞれにファイルを分割するためファイル管理が煩雑であるとともに、どのWebサーバにアクセスするかという制御が必要であり、このため複雑な管理が必要にな

る。また、あるWebサーバに蓄積されているファイルにアクセスが偏ると、そのWebサーバの負荷が増えてしまうという問題もあった。

【0006】また、共有メモリマルチプロセッサ方式においては、CPUのキャッシュメモリにある程度のヒット率が見込まれることを前提に、それに合わせたバスおよびメモリのバンド幅を用意する。しかしながら、共有メモリマルチプロセッサをWebサーバに使用した場合には、CPUのキャッシュにはほとんどヒットしないのでこれがボトルネックとなり、高価で高速なCPUの能力は有効に活用することができなかった。

【0007】さらに、フロントエンド+ファイルサーバ方式は、フロントエンドでHTTPを実行し、必要なファイルをバックエンドネットワークを経由してファイルサーバに要求するというものである("Application of NFS Servers to Strategic Internet/ Intranet Website Design" Technical Report 13, Version 1.0, July 1996, Auspex Systems, Inc. Santa Clara, California)。

【0008】しかしながら、この方式では、フロントエンドとファイルサーバの間の通信がバックエンドネットワークに委ねられている。そのため、バックエンドネットワークは、一般にNFSなど標準の通信規約により行われ、この通信を行うためにフロントエンド、ファイルサーバともにCPUの負荷が増大することが問題となる。

【0009】また、機能分散マルチプロセッサ方式は、ネットワークコントローラ、ファイルコントローラおよびストレージコントローラが、共有バス上でシステムメモリを共有する構造を有しているものである(USP 5, 355, 453)。

【0010】しかしながら、この方式では、複数のネットワークコントローラの読み出しが、全てシステムメモリに集中するので、共有バスおよびシステムメモリのボトルネックになりやすいという問題があった。

【0011】また、疎結合マルチプロセッサ方式は、ハイパーキューブ等で結合されたプロセッサ間結合網を持った疎結合のマルチプロセッサをWebサーバとして使用する("A Scalable and Highly Available Web Server" Proceedings of the IEEE Computer Conference (COMPCON), Santa Clara, March, 1996.)。

【0012】しかしながら、この方式では、プロセッサ間結合網のコストが非常に高いことからシステム全体の価格を押し上げるという問題があった。ところで、TCP/IPは、ネットワークに広く使用されているプロトコルであるが、TCP/IPのパケットをネットワークに送り出す場合は、パケットに対してパリティを計算し、ヘッダ部に付加する必要がある。しかるに、従来は、パリティの計算はCPUによって送り出すデータを読み出して計算を行っていたが、この計算はメモリバン

ド幅にもCPUにも負荷をかけるので望ましくない。

【0013】また、複数のキャッシュ階層の管理においては、第1のキャッシュメモリの内容を第2のキャッシュメモリが全て包含するというマルチレベルインクルージョンという性質が満たされているのが一般的であった。この方法は、第1のキャッシュメモリの速度より第2のキャッシュメモリの速度が遅く、第2のキャッシュメモリはビット単価が安いので容量が大きいという前提において効率的である。

10 【0014】しかしながら、第1のキャッシュメモリ手段、第2のキャッシュメモリ手段の速度の差は大きくないような構成において、キャッシュミスした場合、ディスクアクセスにより多大な遅延時間をとらなう。この場合、キャッシュの実効的な容量を増大させるのが重要であるにもかかわらず、従来のキャッシュ管理方式においては、第1のキャッシュ手段と、第2のキャッシュ手段とに同じ内容が入っていることでメモリ容量にむだが生じることになっていた。

【0015】

20 【発明が解決しようとする課題】上記従来のサーバの構成では、その管理が煩雑であったり、コストパフォーマンスが低いという問題があった。また、パリティの計算はCPUによって送り出すデータを読み出して計算を行っていたが、このような計算はメモリバンド幅にもCPUにも負荷をかけるので望ましくないかった。

【0016】さらに、従来のキャッシュ管理方式においては、第1のキャッシュ手段と、第2のキャッシュ手段とに同じ内容が入っていることでメモリ容量にむだが生じることになっていた。

30 【0017】そこで、本発明は、管理が簡単でかつコストパフォーマンスに優れたサーバを提供することを目的とする。また、本発明は、効率のよいネットワークへの送出手段を提供することを目的とする。さらに、本発明は、階層キャッシュシステムにおける効率のよい記憶領域管理方式をも提供することを目的とする。

【0018】

【課題を解決するための手段】請求項1に係る発明は、少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置(ネットワークインターフェースプロセッサ)と、前記複数の第1の装置と転送手段(内部バス)を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置(ディスクインターフェースプロセッサ)と、前記複数の第2の装置の各々に接続された第3の記憶手段(ストレージデバイス)とを具備し、前記第1の装置に対して与えられるネットワークからの要求に基づく処理を、所定の条件に応じて、前記第1の記憶手段、前記第2の記憶手段または前記第3の記憶手段のうちのいずれかの記憶手段に対して行うことを特徴とするサーバ装置である。

50 【0019】また、請求項2に係る発明は、ネットワー

5

クから要求を与えられる前記第1の装置は、該要求に関連する所定の情報に基づいて、前記第2の装置を選択することを特徴とする。

【0020】さらに、請求項3に係る発明は、前記第3の記憶手段から前記第1の装置に所定のデータが送出される際に、所定のパリティ計算を行い、該得られた計算結果を該所定のデータとともに、前記第1の記憶手段に記憶することを特徴とする。

【0021】また、請求項4に係る発明は、少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備し、前記第1の記憶手段に第1のデータが記憶され、前記第2の記憶手段に該第1のデータに対応する第2のデータが記憶されている場合に、前記第1の記憶手段から該第1のデータが追い出される（実際に追い出す場合のほか、データへのアクセスの無効化等を含む。）ときは、前記第2の記憶手段に記憶された第2のデータの追い出し順位を他のデータの追い出し順位よりも低く設定するための処理を行うことを特徴とするサーバ装置である。

【0022】また、請求項5に係る発明は、少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備するサーバ装置におけるファイル管理方法であって、前記第1の装置に対して与えられるネットワークからの要求に基づく処理を、所定の条件に応じて、前記第1の記憶手段、前記第2の記憶手段または前記第3の記憶手段のうちのいずれかの記憶手段に対して行うことを特徴とするファイル管理方法である。

【0023】さらに、請求項6に係る発明は、少なくとも第1のプロセッサ及び第1の記憶手段を有する複数の第1の装置と、前記複数の第1の装置と転送手段を介して接続される、少なくとも第2のプロセッサ及び第2の記憶手段を有する複数の第2の装置と、前記複数の第2の装置の各々に接続された第3の記憶手段とを具備するサーバ装置におけるファイル管理方法であって、前記第1の記憶手段に第1のデータが記憶され、前記第2の記憶手段に該第1のデータに対応する第2のデータが記憶されている場合に、前記第1の記憶手段から該第1のデータが追い出されるときは、前記第2の記憶手段に記憶された第2のデータの追い出し順位を他のデータの追い出し順位よりも低く設定するための処理を行うことを特徴とするファイル管理方法である。

【0024】本発明によれば、複数用意されたNIP間で通信プロトコルの処理を分散して行うことにより高い

6

性能を得ることができる。また、NIPに用意された第1のバッファメモリ手段をキャッシュとして使用することによりサーバの内部バスのトラフィックを減少させることができ、スケーラビリティが得られるとともに通信のレーテンシを減少させる効果がある。DIPを複数用意することにより、DIPの負荷を分散し、ファイルシステムの処理能力を高めている。また、DIPに用意された第2のバッファメモリ手段がキャッシュとして働くことによりストレージデバイスへの負荷を低減している。

【0025】また、本発明のパリティ計算手段を用いればCPUに負荷をかけずにネットワークにデータを送り出すことが可能となる。さらに、本発明のキャッシュメモリ管理方式を用いれば、キャッシュメモリをより効率的に使えるので、少ないメモリ容量でパフォーマンスの増加が見込まれる。以上のように、本発明によれば、コストパフォーマンスがよく、スケーラブルなWeb等のサーバを構成することができる。

【0026】

20 【発明の実施の形態】以下、図面を参照しながら本発明の実施例を説明する。図1(a)は、本発明の実施形態に係るサーバ装置を用いたシステムの概略構成を示す図である。同図に示すように、本発明に係るサーバ装置1は、ネットワーク2を介してホストマシン3と接続されている。なお、同図において、ホストマシン3は1つのみ接続されているが、特にこれにこだわる必要はなく、複数存在していてもよい。

30 【0027】サーバ装置1は、後述するように、ファイルから読み出したデータをネットワークに送出する機能が実現される。また、サーバ装置1は、ホストプロセッサ（図示せず）が設けられており、当該ホストプロセッサが処理すべきものは、当該ホストプロセッサによって処理される。例えば、WebサーバにおけるCGIの実行は、ホストプロセッサが処理する。

40 【0028】本実施形態においては、ホストプロセッサとサーバ装置との関係を明確にするために、CGI実行を他のファイルからのアクセスと区別する方法について説明する。一般に、CGI実行は、サーバ装置の所定のディレクトリ下で行われる。CGI実行が行われるディレクトリを予めサーバ装置に設定しておき、このディレクトリの下にあるアクセスに対しては、ホストプロセッサにリダイレクトされるように設定しておく。このような設定は、例えば、環境ファイル等により行うことができる。

50 【0029】図1(b)は、本発明の実施形態に係るサーバ装置の構成を示す図である。同図に示すように、サーバ装置1は、ネットワーク2に接続されるネットワークインターフェースプロセッサ21（NIP：Network Interface Processor、以下「NIP」ということもある。）が複数設けられている。ネットワークは、例え

ば、Ethernetが考えられるが、ATM等、他のネットワークでもかまわない。

【0030】各NIP21は、サーバ装置1の内部伝送路（例えば、いわゆる内部バス）であるパラレルリンク12に接続されている。この内部バスは、パラレルリンクに限るものではなく、例えば、シリアルリンク、スイッチ等があげられる。また、このパラレルリンク12には、ディスクインターフェースプロセッサ13（DIP: Disk Interface Processor、以下「DIP」ということもある。）が少なくとも1つ以上接続されており、このDIPにはストレージデバイス14が設けられている。ストレージデバイス14は、例えば、磁気ディスクドライブやディスクドライブ等により実現される。

【0031】図2は、ネットワークインターフェースの構成を示す図である。同図に示すように、NIP11は、プロセッサ21、NIPローカルメモリ23、コード用のROM23、ネットワークコントローラ24およびパラレルリンクインタフェース25から構成される。

【0032】プロセッサ21が実行するプログラムコードは、あらかじめROMに書込まれている。一般に、ROMのアクセス速度は低速であるので、ROMの内容はブートストラッププログラムによりNIPローカルメモリ23にコピーされる。

【0033】NIP11内の上記各ユニットは、例えば、PCIバス26等の標準バスにより接続されている。なお、NIPの構成は、同図に限られるものではなく、プロセッサ21、NIPローカルメモリ22およびPCIバスがブリッジチップ（図示せず）と呼ばれるASICにより接続されるようにしてもよい。NIPのブリッジチップはプロセッサ内のキャッシュメモリ（NIPローカルメモリ内のキャッシュとは異なる。）をスヌープする機能を持つ。スヌープ機構によってDMA転送がプロセッサのキャッシュメモリを更新し、一貫性を保つことができる。

【0034】次に、NIP11の動作について、さらに詳細に説明する。ネットワークからのサーバ装置1に対する要求は、NIP11に伝えられる。プロトコルは、TCP/IPおよびHTTPが使用されるのが一般的であるがこれにこだわるものではない。この要求は、NIP11内のプロセッサ21が所定のプログラム命令を実行することによって、解釈される。

【0035】NIP11内のプロセッサ21は、この要求がホストプロセッサにより実行されるべきものか否かを判断し、ホストプロセッサにより実行されるものであると判断した場合には、該要求をホストプロセッサにリダイレクトする。

【0036】一方、サーバ装置1が実行すべき要求だと判断した場合には、要求されたデータがNIP11内のキャッシュメモリに存在するか否かの判断を行う。本実施形態においては、ローカルメモリ22がキャッシュメ

モリとして機能するが、例えば、プロセッサ21内に設けられる1次キャッシュメモリを含むものであってもよい。

【0037】ここで、所定のデータがキャッシュメモリ内に存在するか否か、すなわち、キャッシュメモリの検索について説明する。図3は、キャッシュエントリの構造を示す図である。同図に示すように、ローカルメモリ内に展開されているキャッシュの内容は、タグ部31とデータ部32とからなる各エントリが、LRUポインタ311によってリンクされており、LRULAST312が最も古くにアクセスされたエントリを、LRUTO313は最も新しくアクセスされたエントリを示している。このリンク構造は、アクセスがあるごとに更新される。ファイル名フィールド312は、エントリが保持しているデータ（ファイル）のファイル名を保持し、タグとして用いられる。TTL（Time To Live）フィールド315は、キャッシュされたデータの有効時間を示す。すなわち、もし、TTLフィールド315の値（TTL値）が現在の時刻と比較して有効期限がきれていると判断される場合には、たとえキャッシュにヒットしてもそのキャッシュの内容は有効でないとみなされる。このTTLフィールド315は、ファイルの属性として付加されており、DIP13からNIPに転送される際に、TTL値も同時に転送される。データ長フィールド316は、エントリのデータの長さを示し、データへのポインタフィールド317は、データの先頭アドレスを示す。

【0038】データ部32は、一定の長さの領域（例えば4Kbyte）に分割されたデータがポインタでリンクされた構造をしており、データ領域はフリーリストで管理される。これは、キャッシュされるデータは可変長なので、データ領域の管理を容易にするためである。

【0039】以上のように、ローカルメモリ22内に展開されたキャッシュの検索は、LRULAST312から順にLRUポインタ311をたどりながらLRUと比較していくことによって行われる。検索を高速化する手法としては、ハッシング等が知られているが、これに限るものではない。

【0040】なお、キャッシュメモリの容量には限界があるので、キャッシュに新しい内容をいれるときにはキャッシュ内の内容を追い出す必要がある。これにはLRUアルゴリズム等を使用して、追い出すべきエントリを決定することができる。追い出すエントリが決定されると、NIP11はDIP13に対して、どのエントリが追い出されるかを通知する。以上がキャッシュの検索の説明である。

【0041】プロセッサ21が、要求されたデータをローカルメモリ22のキャッシュから転送可能であると判断した場合には、プロセッサ21は、ネットワーク送出に必要なヘッダ等の情報を構成し、ローカルメモリ22

10

20

30

40

50

内に書き込む。さらにプロセッサ21は、ネットワークコントローラ24のDMA（図示せず）をセットしてローカルメモリ22からネットワークコントローラ24に転送し、ネットワークコントローラ24が当該情報をネットワークに送出する。

【0042】構成されるヘッダ等の情報は、Webサーバの場合、HTTPプロトコルのヘッダ、TCPのヘッダおよびIPのヘッダが含まれる。ネットワークコントローラは、例えば、DEC社のDEC21140等があげられる。このチップは、PCIインタフェースとDMAコントローラを備えている。

【0043】プロセッサ21が、要求されたデータをローカルメモリ22から供給できないと判断した場合には、プロセッサ21は、DIP13に対して要求を行う。要求は、内部バス（パラレルリンク）12を介して電送される。要求は、必要なファイル情報の他、DIP13からNIP11にデータを転送する際のNIP11側のバッファメモリのアドレスが含まれる。

【0044】NIP11のプロセッサ21は、該要求を複数のDIPのうちのどのDIPに送出するかを決定する。本実施形態においては、ファイル名に基づいてハッシュ値を求め、該ハッシュ値に基づいて決定される。例えば、DIP13が3個設けられている場合、ハッシュ関数としてファイル名のキャラクタコードの和のmod 3をとったものが使用される。NIP11からDIP13への要求の他、上述したようななどのエントリーを追い出したかという情報をNIP11からDIP13に通知する場合にも、同様の方法でDIP13を決定する。

【0045】図4は、パラレルリンクインタフェース25の構成を示す図である。同図に示すように、パラレルリンクインターフェース25は、パラレルリンク制御部41、パリティ計算部42、Receiveパリティバッファ43、Receiveバッファ44、Sendバッファ45、ReceiveDMA46、SendDMA47およびPCIインタフェース48より構成される。

【0046】パラレルリンク制御部41は、DIPからの通信のうち、宛先が該当するNIPに対する通信（メッセージ）のみを通過させ、パリティ計算部42に送出する。パリティ計算部42は、パリティの計算を行い、一定のワード数ごとにその結果をReceiveパリティバッファ43に書き込む。Receiveパリティバッファ43は、例えば、256ワードに1回の書き込みでは256ワード転送毎に書き込まれる。ここでのワードは、TCPのパリティの形式にあわせて16ビットで1ワードとする。

【0047】ReceiveDMA46は、Receiveパリティバッファ43に書込まれたデータを、NIPローカルメモリ22に転送する。TCPパケットをNIP11のプロセッサ21が作成する場合、Receive

パリティバッファ43から転送されたパリティデータを使用する。例えば、TCPパケットが1024ワードのとき、その1024ワードに対応する4ワードのパリティデータからパケットのパリティを作成する。1024ワードがパリティデータの256ワードのアライメントにあわない場合には、アライメントにあわない端数部分のパリティを計算する必要がある。Receiveパリティバッファ43を設けることによって、NIP11のプロセッサ21によるパリティ計算の負荷を著しく減少することができる。また、256ワード毎にパリティを付加することにより、TCPパケットのサイズに柔軟に対応することができる。

【0048】DIPに対する要求は、NIP11のプロセッサ21によりコマンドが構成され、SendDMA47によりSendバッファ45に転送され、さらにパラレルリンク12に送出される。

【0049】図5は、NIP11にネットワークから要求がきた場合の処理を説明するための図である。同図において、まず、ネットワークコントローラ24を介してネットワークから要求を受け付けると（STEP51）、CGIか否かが判断される（STEP52）。要求がCGIであると判断された場合には、NIP11のプロセッサ21は、CGIを実行するように設定されていないため、ホストマシンに処理を要求する（STEP53）。一方、CGIでないと判断された場合には、次に、NFSであるか否かが判断される（STEP54）。

【0050】NFSであると判断された場合には、要求を送出するDIPを決定し（STEP55）、該決定されたDIPにNFSを要求する（STEP56）。一方、NFSでないと判断された場合には、上述のキャッシュの検索を行う（STEP57）。キャッシュを検索することによりヒットした場合には、TTLをチェックし（STEP58、59）、有効期限内か否かを判断する（STEP510）。TTL値が有効期限内であると判断された場合には、NIP11のNIPローカルメモリ22からデータをネットワークに送出する（STEP511）。一方、有効期限内でないと判断された場合には、要求を送るDIPを決定し（STEP512）、該決定されたDIPに要求を送出する。

【0051】図6は、NIP11がDIP13からの応答を受け付けた場合の処理を説明するための図である。同図において、DIPから応答を受け付けると（STEP61）、NFSからの応答であるかが判断される（STEP62）。NFSからの応答であると判断された場合には、そのままネットワークに応答を返す（STEP63）。NFSからの応答でない場合に、要求したデータが帰ってきた場合には、NIP11は、NIPローカルメモリ22内のキャッシュがいっぱいであるか否かを判断する（STEP64）。キャッシュがいっぱい

である場合には、追い出すべきエントリを決定し、該エントリを追い出すとともに（STEP 65）、追い出したエントリをDIPに通知する（STEP 66）。キャッシュがいっぱいでないと判断された場合、またはエントリの追い出し処理を行った後、NIP11は、キャッシュに該送られてきたデータをキャッシュに書き込むとともに、ネットワークに送出する（STEP 67）。

【0052】次に、DIP13の詳細について説明する。図7は、DIP13の構成を示す図である。同図に示すように、DIP13は、プロセッサ71、DIPローカルメモリ72、コード用ROM73、SCSIインタフェース74、パラレルリンクインタフェース75およびファイルバッファメモリ（図示せず）から構成される。

【0053】DIP13内の上記各ユニットは、PCIバス等の標準バスにより接続されている。なお、図7に示す接続方法のほか、プロセッサ71とDIPローカルメモリ73およびPCIバス11を接続するために、ブリッジチップと呼ばれるASICを用いてもよい。

【0054】図8は、DIP13のパラレルリンクインタフェース75の構成を示す図である。同図に示すように、パラレルリンクインターフェース75は、パラレルリンク制御部81、Receiveバッファ82、Sendバッファ83、ReceiveDMA84、SendDMA85およびPCIインタフェース86より構成される。

【0055】パラレルリンク制御部41は、NIP11からの通信のうち、宛先が該当するDIPに対する通信（メッセージ）のみを通過させ、Receiveバッファ82に送出する。

【0056】ReceiveDMA84は、Receiveバッファ42に書込まれたデータを、PCIインタフェース86を介して、DIPローカルメモリ72に転送する。

【0057】NIP11に対する応答は、DIP12のプロセッサ71によりコマンドが構成され、SendDMA85によりSendバッファ83に転送され、さらにパラレルリンク81に送出される。

【0058】DIP13は、パラレルリンクインタフェース75を介して、NIP11からのデータの要求またはキャッシュエントリの追い出し通知を受け取る。DIP13は、NIP11からのデータ要求を受け取ると、DIPローカルメモリ72に該要求されたデータが存在するか否かを調べる。DIPローカルメモリにおいても、図3に示したNIPローカルメモリと同様に、キャッシュとして機能させるための各種フィールドを有する。

【0059】図9は、DIPの動作を説明するための図である。同図において、DIP13は、いずれかのNIP11から通知を受け取ると（STEP 91）、該通知

が追い出し通知であるか否かを判断する（STEP 92）。該通知が追い出し通知であると判断された場合には、キャッシュに該当するエントリを検索する（STEP 93）。検索の結果、該当するエントリがあるか否かを判断し（STEP 94）、ないと判断された場合には、特になにも行わない。一方、該当するエントリがあると判断された場合には、該当するエントリをLRUTOPにつなぎ直す（STEP 95）。これは、DIPローカルメモリ72のキャッシュのエントリを追い出す場合に、NIPローカルメモリのキャッシュから追い出されたエントリであるか否かを考慮するものである。

【0060】すなわち、NIPローカルメモリのあるエントリが追い出し通知を受けた場合に、そのエントリのコピーがDIPローカルメモリに存在する場合には、DIPローカルメモリの対応するエントリの追い出しの優先順序を低く、換言すれば、追い出されにくくなるように設定する。これによって、NIPローカルメモリにあるエントリとDIPローカルメモリにあるエントリの重複が起こる確率を少なくすることができる。なお、どのエントリを追い出すかを決定するには一般にLRUアルゴリズムが使用される。

【0061】STEP 92において、追い出し通知でないと判断された場合、NFS要求であるか否かが判断され（STEP 96）、NFS要求である場合にはNFS（Network File System）処理を行う（STEP 97）。一方、NFSでないと判断された場合には、DIPローカルメモリのキャッシュを検索し（STEP 98）、ヒットしたか否かを判断する（STEP 99）。ここで、ヒットした場合には、TTLのチェックを行い（STEP 910）、有効期限内であれば（STEP 911）、キャッシュからデータを読み出して、要求もとのNIPに応答する（STEP 912）。

【0062】キャッシュの検索がヒットしなかった場合、または有効期限内でないと判断された場合には、キャッシュから追い出すエントリを決定し、キャッシュエントリを空け（STEP 913）、ストレージデバイス14より読み出して（STEP 914）、該空いたエントリに読み出されたデータを書き込むとともに、要求もとのNIPにデータを送出する（STEP 915）。

【0063】このように、ファイル（データ）をネットワークに送出する場合は、NIP11およびDIP13が協調して動作し、第1のバッファメモリ手段であるNIPローカルメモリ22および第2のバッファメモリ手段であるDIPローカルメモリ72をそれぞれキャッシュとして利用し、該キャッシュにデータが存在しない場合には、ストレージデバイス14からデータを読み出している。なお、ストレージデバイス14からの読み出しは、DIPのファイルシステム（例えば、プロセッサ71によるプログラムの実行により実現される）によって行われる。

10

20

30

40

50

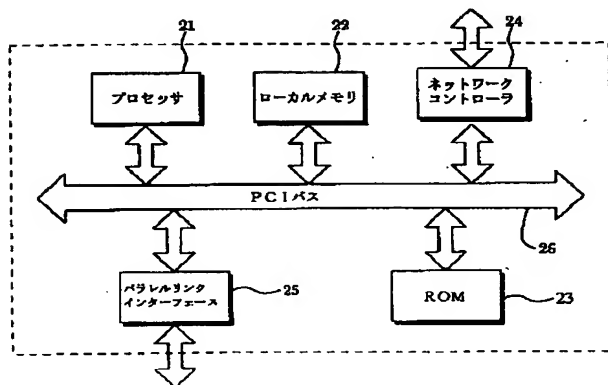
【0064】次に、本発明に係るサーバへの書き込み動作について説明する。書き込みは、ネットワークから行われる。ネットワークからの書き込み要求は、NIP11が適切なDIP13に要求を転送する。適切なDIPを選択するのはNIP11が上述したファイル名によるハッシングを利用して行われる。NIPから要求を与えられたDIP13は、ストレージデバイスに対して要求されたデータの書き込みを行う。また、ネットワークからDIP13に接続されたストレージデバイス14にはNFSで読み出し、書き込みができるようにしておく。この機能を実現するために、NIP11は、ネットワークからのNFS通信をDIP13に中継する処理が行われる。

【0065】ネットワークからNFSでDIP13に接続されたストレージデバイス14に読み書きができることで、ネットワーク上のホストマシンから見てサーバ装置のストレージデバイス14がリモートファイルシステムとして実現される。ホストマシンからリモートファイルシステムとしてみえることで、ホストマシン上にある従来のソフトウェアを有効に活用することが可能である。本実施例ではNFSをファイル書き込みのプロトコルとして説明したが、他のネットワークプロトコルを使用してもよい。

【0066】

【発明の効果】以上説明したように、本発明によれば、ストレージデバイスに蓄積されたファイルを効率よくネットワークに送り出すことができる。また、階層化されたキャッシュ管理方式により、高いヒット率を得ること

【図2】



ができる。さらに、データをネットワークに送出するとき必要になるパリティ計算をデータの転送中に行うことでプロセッサの負荷を減らすことができる。

【図面の簡単な説明】

【図1】 本発明の実施形態に係るシステムの概略構成およびサーバ装置の構成を示す図。

【図2】 ネットワークインターフェースの構成を示す図。

【図3】 キャッシュエントリの構造を示す図。

10 【図4】 NIPの平行リンクインタフェースの構成を示す図。

【図5】 NIPにネットワークから要求がきた場合の処理を説明するための図。

【図6】 NIPがDIPからの応答を受け付けた場合の処理を説明するための図。

【図7】 DIPの構成を示す図。

【図8】 DIPの平行リンクインタフェースの構成を示す図。

【図9】 DIPの動作を説明するための図。

20 【符号の説明】

1…サーバ装置

2…ネットワーク

3…ホストマシン

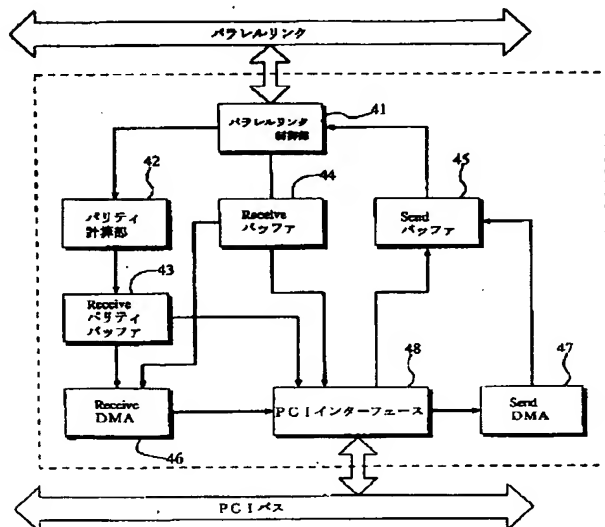
11…ネットワークインターフェースプロセッサ (NIP)

12…平行リンク (バス)

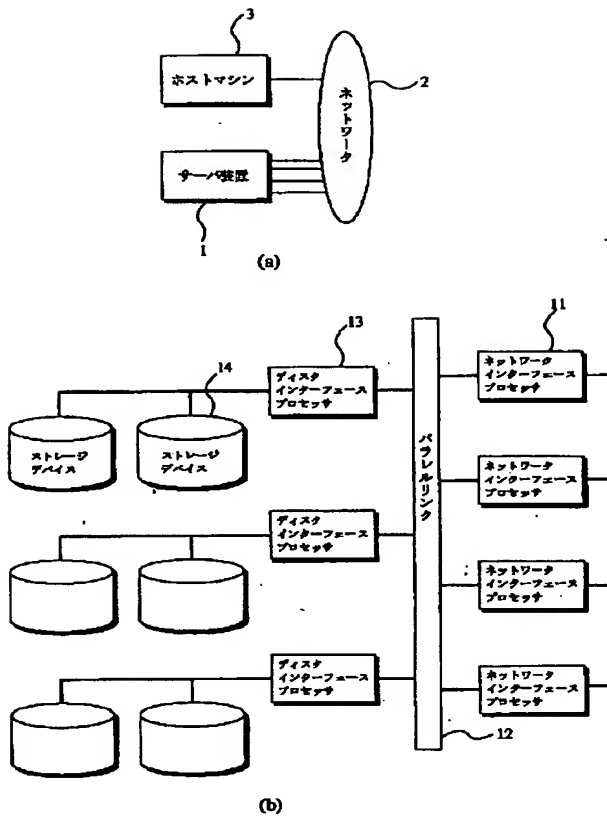
13…ディスクインターフェースプロセッサ (DIP)

14…ストレージデバイス

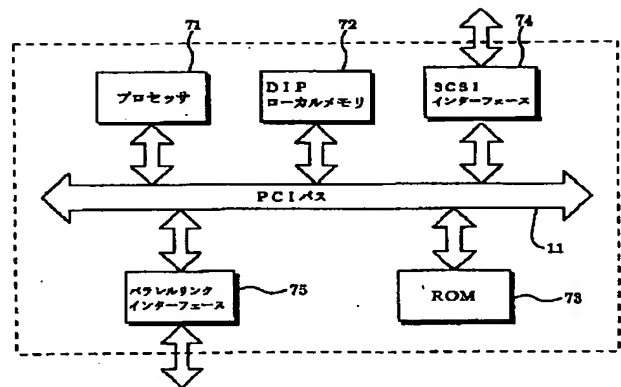
【図4】



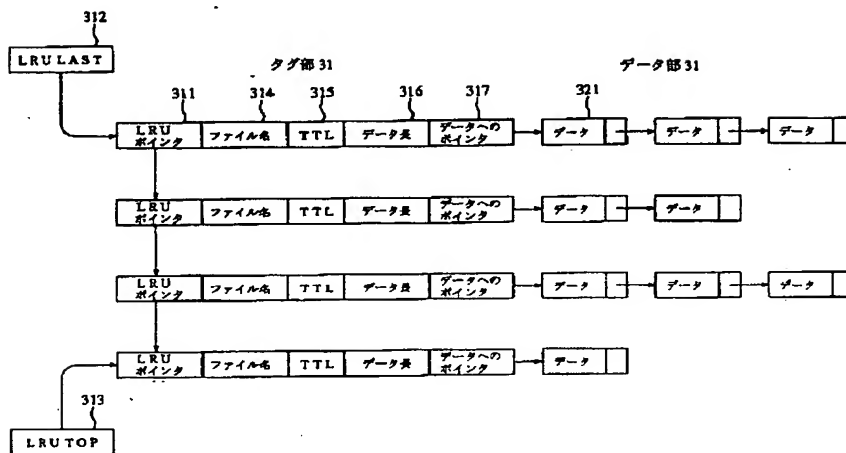
【図1】



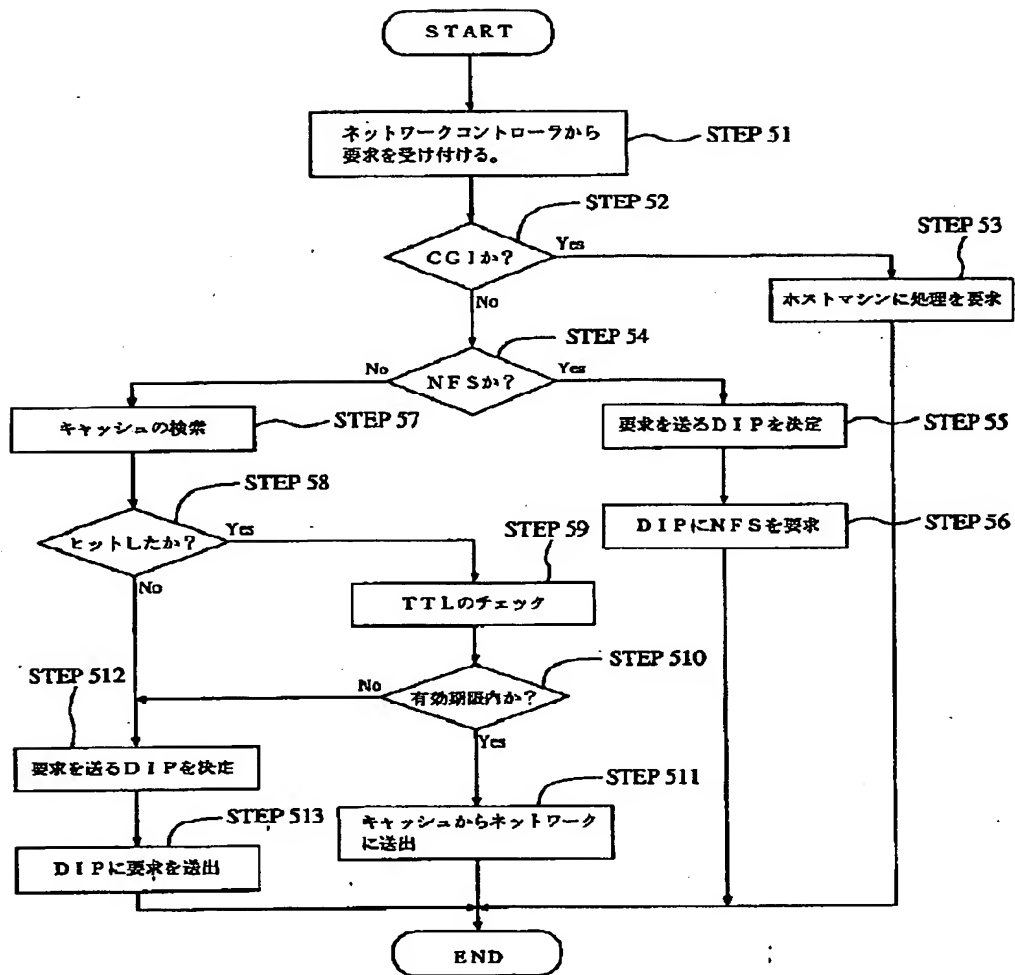
【図7】



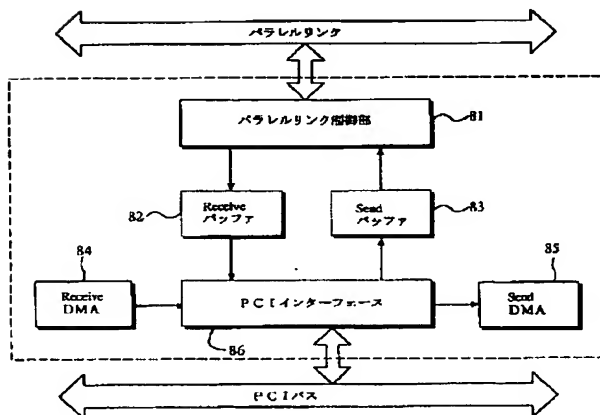
【図3】



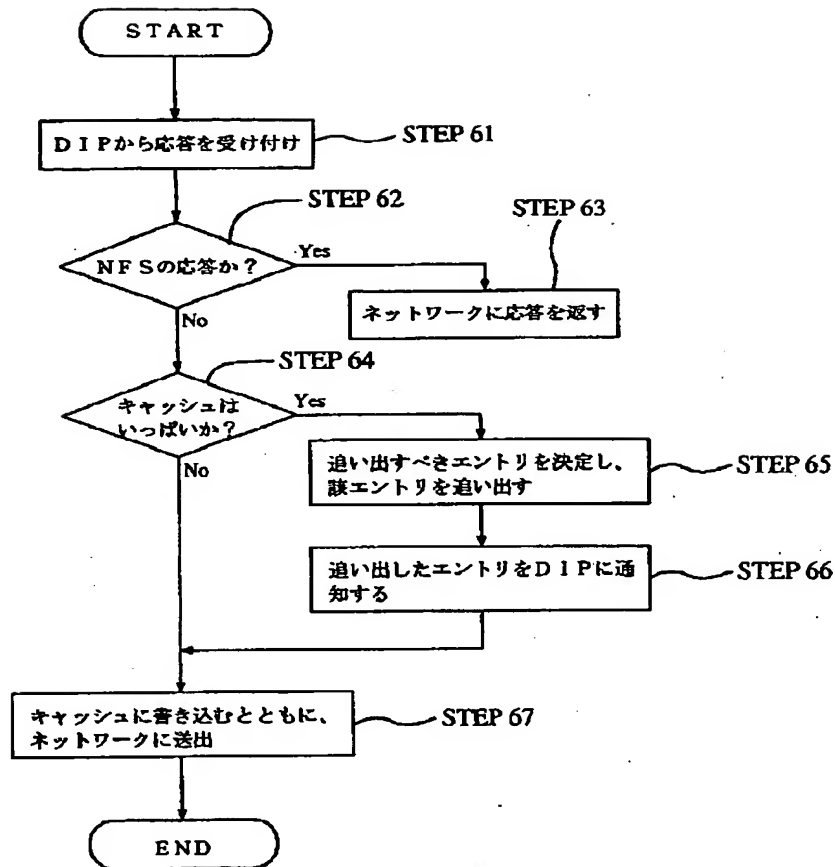
【図5】



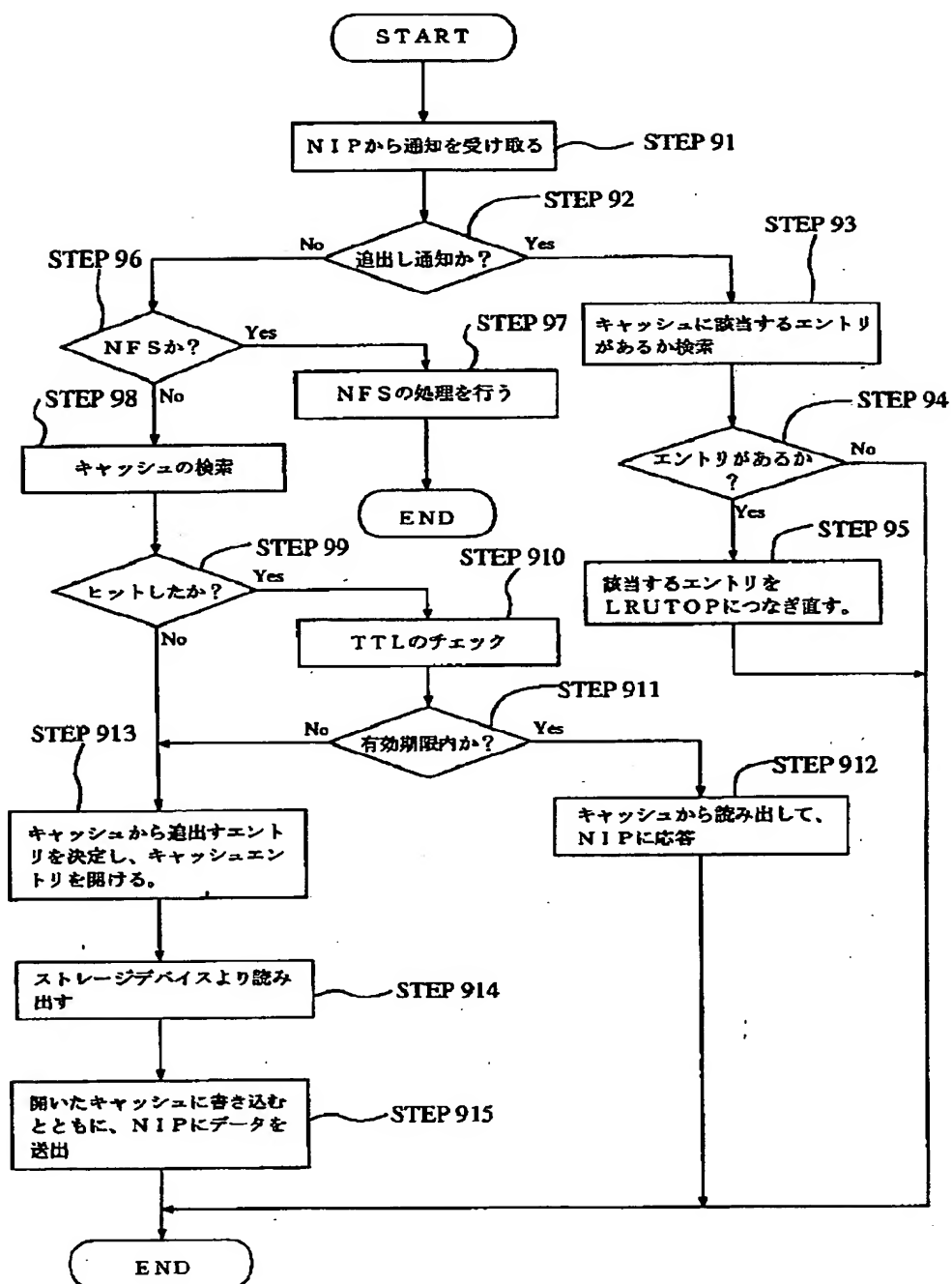
【図8】



【図6】



【図9】



フロントページの続き

(72)発明者 前田 誠司

神奈川県川崎市幸区小向東芝町1番地 株
 式会社東芝研究開発センター内